of params we need to estimate (and helps lessen curse of dimensionality).

# Factor Model:

- Assume $E[\vec{X}] = 0$
- Assume there is a $q$-dimensional random vector $\vec{F}$ with $E[\vec{F}] = 0$, $Var[\vec{F}] = I$, $q < p$
- $\vec{F}$ is latent/unobserved/hidden, $\vec{X}$ is manifest/observable
- Assume $\underset{p \times 1}{\vec{X}} = \underset{p \times q}{W} \underset{q \times 1}{\vec{F}} + \underset{p \times 1}{\vec{\varepsilon}}$, with $E[\vec{\varepsilon}] = 0$, $Var(\vec{\varepsilon}) = \Psi$ ($\Psi = $ diag matrix), $Cov(\vec{\varepsilon}, \vec{F}) = 0$

So under these assumptions

$$Cov(X_i, X_j) = Cov\left(\sum_{k=1}^{q} W_{ik} F_k + \varepsilon_i, \sum_{\ell=1}^{q} W_{j\ell} F_\ell + \varepsilon_j\right)$$

$$= Cov\left(\sum_{k=1}^{q} W_{ik} F_k, \sum_{\ell=1}^{q} W_{j\ell} F_\ell\right)$$

$$= \sum_{k=1}^{q} \sum_{\ell=1}^{q} W_{ik} W_{j\ell} Cov(F_k, F_\ell)$$

$$= \sum_{k=1}^{q} W_{ik} W_{jk} \quad \left(\text{since } Cov(F_k, F_\ell) \text{ is } 0 \text{ when } k \neq \ell \text{ and } 1 \text{ when } k = \ell\right)$$

So $X_i$ and $X_j$ are correlated when they load on the same factors

$$\text{Var}(\vec{x}) = \text{Var}(w\vec{F} + \vec{\epsilon}) = \underset{pxq \ qxp}{ww^T} + \underset{pxp}{\psi}$$

## How to Estimate:

If we know $\psi$, then $\text{Var}(\vec{x}) - \psi = ww^T$
where $ww^T$ is symmetric and positive definite
so $w = ud^{1/2}$ and we have

$$\text{Var}(\vec{x}) - \psi = udu^T$$

To estimate $\psi$
• regress each variable on the others and find MSE of estimates (can be initial estimate for $\psi$)

## How to pick number of factors:

$$q = \text{rank}(\text{Var}(\vec{x}) - \psi) = \text{rank}(ww^T)$$

• look at eigenvalues of $\widehat{\text{Var}(\vec{x})}$ and stop when they get small (like with scree plot, and look at elbow)
• test goodness of fit, $||\widehat{\text{Var}(\vec{x})} - (\hat{w}\hat{w}^T + \hat{\psi})||$ is small enough to be noise
• if we assume distribution for $\vec{F}$, we can get a dist for $\vec{x}$. Evaluate accuracy of dist using log-likelihood or something.